

Integrating Natural Language and Gesture in a Robotics Domain*

Dennis Perzanowski, Alan C. Schultz, and William Adams
Navy Center for Applied Research in Artificial Intelligence
Naval Research Laboratory
Washington, DC 20375-5337

ABSTRACT

Human-computer interfaces facilitate communication, assist in the exchange of information, process commands and controls, among many additional interactions. For our work in the robotics domain, we have concentrated on integrating spoken natural language and natural gesture for command and control of a semi-autonomous mobile robot.

We have assumed that both spoken natural language and natural gesture are more user-friendly means of interacting with a mobile robot, and from the human standpoint, such interactions are easier, given that the human is not required to learn additional interactions, but can rely on “natural” ways of communication. So-called “synthetic” methods, such as data gloves, require additional learning; however, this is not the case with natural language and natural gesture. We, therefore, rely on what is natural to both spoken language when it is used in conjunction with natural gestures for giving commands.

Furthermore, we have been integrating these interactions with the robotics components as the robotics system is being developed. The interface is not treated as an ad hoc add-on or patch. By doing so, we believe the interface will be more robust and because it is being integrated during system development, we hope to achieve a more seamless interface, one which both acts and feels as an integral part of the robotics application.

In this paper, we will discuss the kinds of interactions which our system is currently capable of performing. We will also discuss the processing of the various input to produce an appropriate robotic response. And finally, we will discuss what future kinds of interactions we would like to incorporate into the system, and what will be required to achieve these results.

KEYWORDS: *biomimetics, gesture, human-computer interfaces, natural language, robotics*

1. INTRODUCTION

Our research implementing a natural language and gestural interface to a semi-autonomous robot is based on two assumptions. The first, or linguistic, assumption is that certain types of ambiguity in natural language can be resolved when gestures are incorporated in the input.

For example, a sentence such as “Go over there” is devoid of meaning unless it is accompanied by a gesture indicating the place where the speaker wishes the hearer to move. Furthermore, while gestures are an integral part of communication [1], our second, or gestural, assumption is that stylized or symbolic gestures place a heavier burden on the human, frequently requiring a learning period, since such gestures tend to be arbitrary in nature. Natural gestures, i.e. gestures that do not require learning and which any human might produce as a natural co-occurrence to a particular verbal command, are simpler means of imparting certain kinds of information in human-computer interaction. With systems that have fairly robust vision capabilities, natural gestures obviate the need for additional interactive devices, such as computer terminals, touchscreens, or data gloves. So from a linguistic and gestural standpoint, certain utterances, such as those that involve movement or location information, can be disambiguated by means of natural, accompanying gesture [2].

Furthermore, ample evidence from related research [3] indicates that there is a close relationship between speech and gesture. A natural language and gestural interface, therefore, should utilize and maximize this known relationship.

For this study, we limit ourselves to two types of commands: commands that involve direction, e.g. “Turn left,” and those that involve locomotion, e.g. “Go over there.” For such commands, environmental conditions permitting, people communicate with each other by pointing to objects in their surroundings, or gesturing in the specified direction. Granted, if the environment or meteorological conditions are not favorable, as for example when it is too dark to see or if foggy or heavy precipitation prevails, humans may rely on other methods to communicate, which will not concern us here. However, given a more or less ideal environment, human to human communication typically involves the use of natural language and gesture, and it is this type of interaction that we have emulated in our human-computer interface to a semi-autonomous robot.

*This work is funded in part by the Naval Research Laboratory and the Office of Naval Research. This publication was prepared by United States Government employees as part of their official duties and is, therefore, a work of the U.S. Government and not subject to copyright.

For the kinds of interaction that we have outlined above, touchscreens or data gloves also allow humans to communicate and talk about so-called “deictic” elements (to be defined shortly) in various computer applications. While such devices may be appropriate in many applications, we have excluded them from our work.

These “synthetic” methods of interaction require additional learning, and may not, therefore, be the easiest or most natural ways of interacting. Ultimately, we are concerned with utilizing and transferring as much of human-human interaction as is possible to human-robot interactions. However, while it is yet to be proven that all, or at least, much of human-human interactions is appropriate in human-robot interactions, we work with the assumption that people find it easier and will tend to react and communicate more naturally when given this opportunity, even when interacting with robots.

For the purposes of this investigation, we have limited ourselves to robotic commands, and particularly those that involve commands of movement and of distance. These linguistic utterances typically contain so-called “deictic” elements. “Deictics” are linguistic strings that refer to objects in the discourse which in turn usually refer to objects in the real world.

For example, in the sentence “the box in the corner is red,” the subject of the sentence “the box in the corner” can be analyzed as a deictic element if one exists in the same environment as the speaker and/or hearer of this utterance. If the intended referent, namely “the box,” does not exist in this environment, either the speaker is playing some sort of linguistic trick, or the utterance is uninterpretable.

More typically, deictic elements are characterized by the presence of such words as “this” or “that” for objects, as in the expressions, “this box is red,” “that is a blue box,” or “here” or “there” for locations, as in “bring it over here,” or “the waypoint over there.” We limit ourselves in this research to statements in which deictic elements exist in the world of the speaker and hearer.

Based on our initial assumptions, we have also chosen to limit our consideration to those natural gestures that involve gross movements of the hands and arms. Granted, someone can disambiguate “Go/Move over there” by moving one’s head in a particular direction, or by moving one’s eyes.

Because our system cannot currently handle such complex and minute movements for the purposes of disambiguation, we limit ourselves to the grosser movements, such as hand and arm movements, which our vision system can discriminate. However, we believe that there are no comparable natural language and gesture robotic interfaces that automatically produce robot controls by combining natural language and natural gesture.

2. THE NATURAL LANGUAGE AND GESTURE INTERFACE

This work is based on research already underway at the Navy Center for Applied Research in Artificial Intelligence (NCARAI) [4]. To process verbal or audial commands issued to a semi-autonomous robot, the natural language interface, depicted in Figure 1, utilizes a 70-word speech vocabulary with an input range of approximately 11,000 utterances.

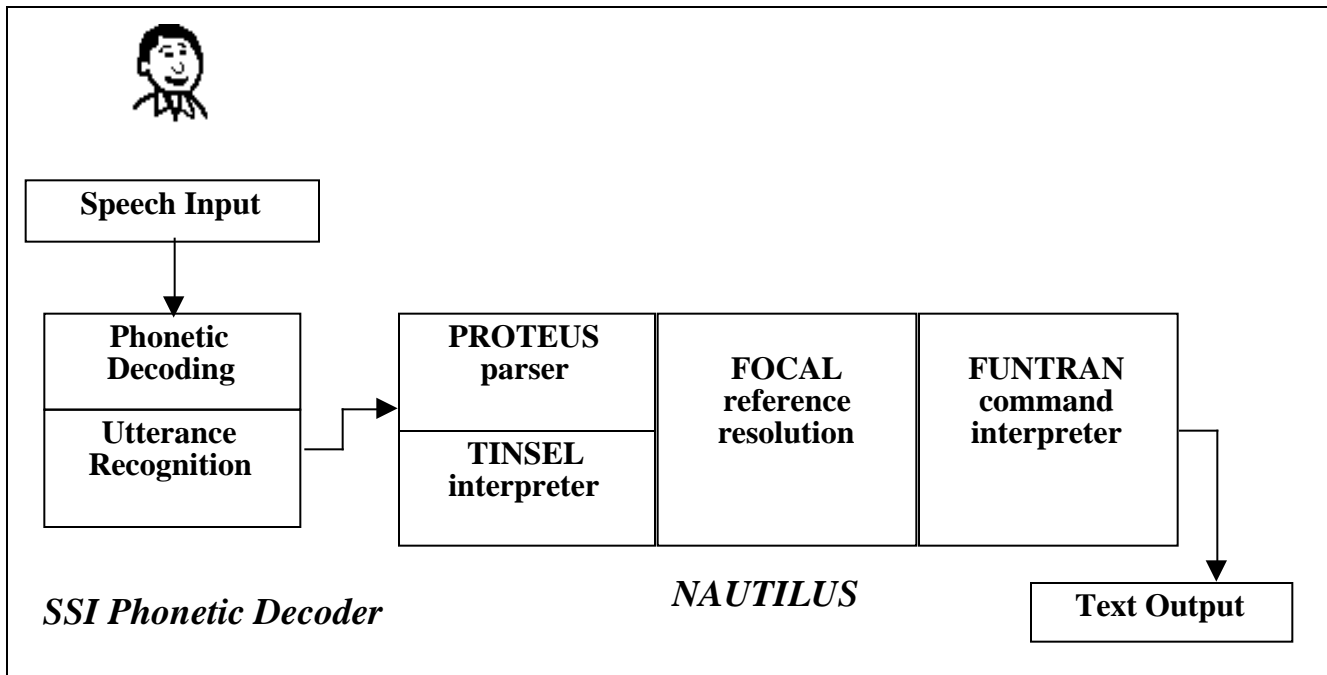


Figure 1. Natural language processing of command utterances

2.1. Natural Language Processing

The auditory signal is converted to a text string by the speech recognition system, a PE 200, manufactured by Speech Systems, Inc. The textual string is then parsed by our natural language processing system, NAUTILUS, developed in-house at NCARAI [5].

We have opted for robust natural language parsing, starting with syntactic analysis of the input string, since one of our project's research objectives [6] has been to port the grammar and parsing mechanism to various applications and domains. Furthermore, as this project develops and tackles more complex commands involving natural language and gesture, we believe it will be necessary to have a full parse which will need to utilize timestamp information in the speech signal, as well as lexical and phrasal information which will further need to be coordinated with timestamp information in the gestural input. We, therefore, believe that robust parsing can provide us with richer interpretations, so that our ultimate desire to understand more complex utterances in conjunction with gestures will be attainable.

In Figure 2, we see how the various linguistic modules analyze the natural language input and map it to a corresponding gestural input. The auditory signal is analyzed and mapped to a syntactic string. During the parsing process, the syntactic string is semantically interpreted in several of the NAUTILUS modules, and the resulting semantic representation is submitted to the command module on the semi-autonomous robot for further processing.

The mapping of the semantic interpretation, the Lisp-like structure on the left-hand side of Figure 2, and the perceived gesture, the numerical string on the right-hand side of Figure 2, are compared to produce some kind of action which is best understood after a consideration of the gesture processing, to which we now turn.

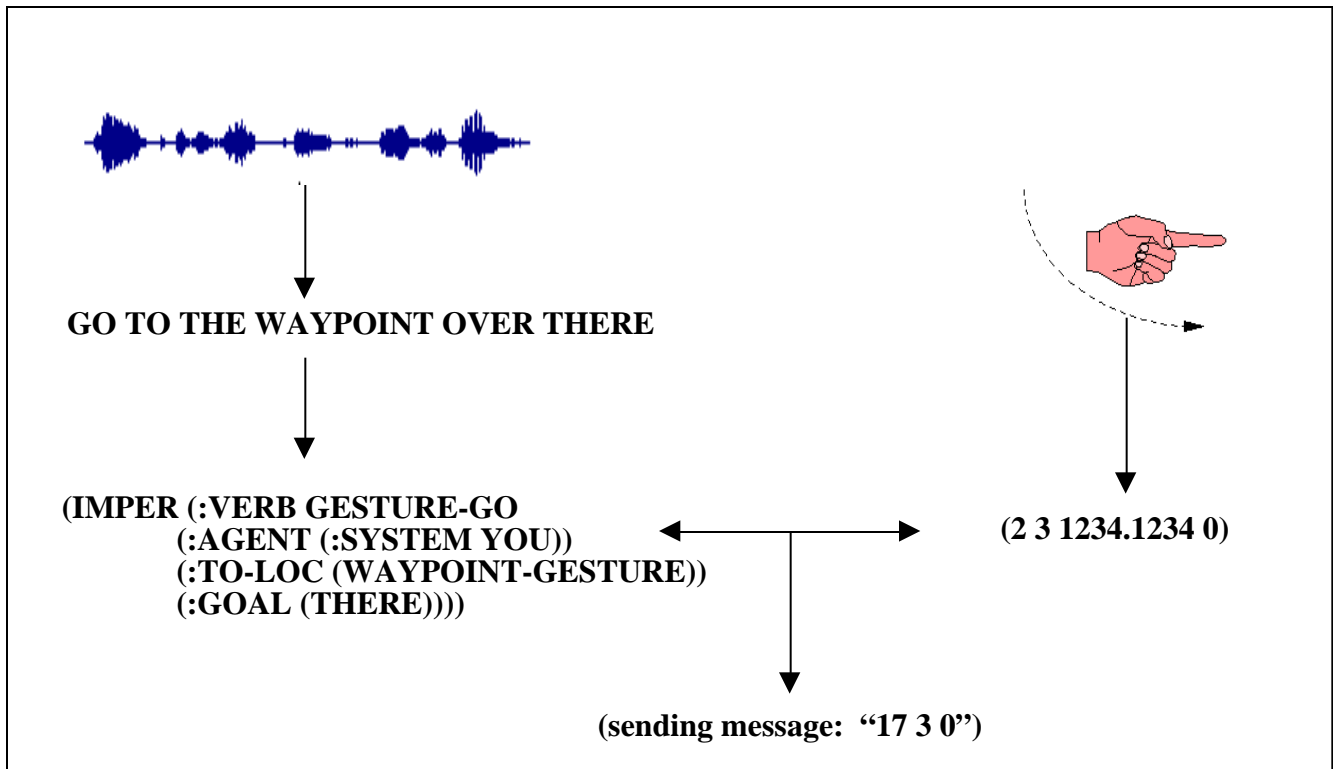


Figure 2. Integration of natural language processing and gesture analysis

2.2. Gesture Processing

The semi-autonomous robot in this investigation, which we have nicknamed "Coyote," is a mobile Nomad 200, manufactured by Nomadic Technologies, Inc., and equipped with 16 Polaroid sonars and 16 active infrared sensors. It is capable of detecting and incorporating gesture into various verbal commands (Figure 3).



Figure 3. A Nomad 200 mobile robot with mounted camera

By means of the top-mounted camera, the robot is capable of sensing vectors and measured line segments that the human might gesture during the various commands. The robot is linked as a UNIX workstation via a radio ethernet connection to the natural language processing modules and to the command modules of the robot.

A process running on the robot is used to determine the gestures given by the human user. The gestures are detected with a structured light rangefinder, which emits a horizontal plane of laser light 30 inches above the floor. A camera mounted above the laser is fitted with a filter tuned to the laser frequency. The camera observes the intersection of the laserlight with any objects in the room, and the bright pixels in the camera's image are mapped to XY coordinates.

Periodically, the data points from one camera frame are used to compute an average distance from the objects seen. The points are then sorted into clusters, and any cluster sufficiently closer than the average and of appropriate size is designated as a hand. Hand locations are stored for multiple frames until no hands are found or a maximum number of frames are used.

The hand locations across the frames are ordered into one or two trajectories. Trajectories are built incrementally by grouping each hand location of each frame with the trajectory with which it best aligns.

Completed trajectories are checked to see if they are in motion or are stationary, and then logically compared to determine if the overall gesture is valid and if so which gesture was made. The valid gestures are queued and when the multimodal software needs to check for a gesture,

it queries the gesture process which returns the most recent gesture from the queue. This is the string, a kind of data structure, on the right-hand side of Figure 2.

2.3. Mapping Speech Input with a Gesture

We now return to the natural language side of processing the input. An utterance, such as "Go/Move over there", is inherently ambiguous without additional information. Cues must be supplied, namely an accompanying gesture for the sentence to be completely understood. The natural language system can parse the utterance and give an adequate semantic interpretation to it, but without a visual cue, such as a hand gesture in a particular direction, the sentence is ultimately meaningless. Given the vision capability on the Nomad 200 robot and the processing as outlined above, when the sensors on the robot detect a vector within the limitations of its light striping sensor, and a verbal command involving movement in some direction is parsed, a query is made of the gesture process on the robot to see if some gesture has been perceived. The two, namely the semantic interpretation of the verbal command and the data structure containing information about the perceived gesture, are then mapped to a message, the final output in Figure 2. This message, consisting of a string of digits, is sent back to the robot, where it is further processed in order to produce an appropriate action. The mapping of the speech input and the perceived gesture is a function of the appropriateness or inappropriateness and the presence or absence of a gesture during the speech input. Thus, appropriate error messages can be produced if the correct mapping of verbal command and gesture is not made. For example, if the verbal command "Go/Move over there" is interpreted, but no gesture is perceived, a DECTalk speech synthesizer, one of the peripheral modules in the system, produces an appropriate audial error message, such as "Where?".

3. ADDITIONAL TYPES OF INPUT

To handle specific objects from the environment in the various commands which the robot is capable of processing, we have introduced these objects into the semantic component of the natural language processing system. Thus, when a sentence such as "Go to waypoint two" is uttered, and whether or not an object is pointed to in the environment, a meaningful utterance is obtained. The knowledge base of the robot is consulted, and given the fact that an appropriate object (waypoint two) exists in the semantics of the natural language component and in the robot's knowledge base, the robot then moves toward the known object in the room. In this case, since a referent in the real world exists and is known, and has been uniquely identified in the speech signal, a gesture is redundant, and can be ignored. If the human utters the sentence "Go to the waypoint over there," after natural language processing successfully parses and interprets the utterance, it queries the robot's knowledge base to check if there is an object of

such a description as a “waypoint” located in that area of the room, assuming a gesture has been made. If the query receives an affirmative response, the robot moves off to the intended goal. If the human points to some location in the room where there is no known waypoint, then an appropriate error response is produced, such as “I don’t understand. There is no waypoint in that direction.” Of course, if no gesture was perceived, the robot responds appropriately that one is required.

Additional commands, such as “Back up/Move forward this far” are handled in a similar fashion. If the camera mounted on the top of the robot senses a measured line segment, the robot moves the gestured distance in the intended direction; however, without an appropriate gesture, such a sentence evokes the error response, “How far?”.

Directional commands are treated in much the same way. For example, the robot can be told to turn in any direction an arbitrary number of degrees. Such utterances as “Turn 30 degrees to the left,” “Turn to your left/right” or “Turn to my left/right” produce appropriate robotic responses. However, if the human issues a contradictory gesture while uttering these commands, the robot responds accordingly, stating that a contradictory gesture was perceived, and no further action is taken at that time. Likewise, if the robot is told to “Turn this/that way,” a gesture must be perceived; otherwise, an error message results, stating that no gesture was observed and one is required with such a command.

As anyone involved with speech recognition systems knows, numbers are still extremely difficult to understand and pose major problems. However, since the focus of our work here has been on integrating natural language and gesture and not totally on speech understanding, we await a more sophisticated and advanced speech recognition system to continue our work.

4. SOME CONSTRAINTS OF THE CURRENT SYSTEM

As we continued to process directional commands involving turning and moving toward objects, we realized that some of the directional capabilities of our system were being severely constrained because of certain physical constraints imposed on us by our speech recognition system.

We have already seen one limitation: namely, the problem that speech understanding systems have with the processing of numbers. In part this problem results from the auditory similarity of such numerical strings as “thirteen” and “thirty.” These strings can sound strikingly similar even to human listeners. This is compounded by the fact that there is very little context to offer clues for processing such strings in such an utterance as “Turn left ____ degrees.” However, we leave these problems for researchers in speech recognition and do not consider it a problem for us to solve in robotics interfaces.

The system that we currently employ, the SSI PE 200, requires that the human user be tethered to the speech recognition device. It requires the user to wear a headset that is directly connected by a cable to the speech recognition system. When any orientation of speaker and robot, other than face-to-face orientation, results, we are forced to maneuver the robot so that its vision system is capable of seeing the user and any gestures produced in a face-to-face orientation. At times, this requires some strategic re-alignment of robot and human user. Also, a command such as “Turn to my left/right” suddenly takes on an entirely different meaning when speaker and hearer are no longer oriented face-to-face. With a speech recognition system that permits user mobility, the user will be able to move more freely in the environment, thereby increasing the types of interaction. However, this will also require more sophisticated auditory sensors on the robot.

With increased user mobility, the orientation of human and robot can vary greatly; consequently, the robot must be aware of where the human is with respect to itself. Granted, it might try to visually scan the area whenever a verbal command is perceived. However, this does not seem the most efficient way to handle this problem.

Humans, for example, do not necessarily visually scan an area to see where a command is being produced. We rely on our stereophonic hearing capabilities. We then interpolate the direction from which the command is being issued and react accordingly. In this way, we can compensate for the new orientation of human and robot whenever intervening movement of either participant in the interaction may change their orientation to each other.

Furthermore, because of the vision capabilities that are currently used by our robot, we are limited to vectoring gestures or gestures that segment a line. A more sophisticated visual system, such as one that is capable of discriminating hand shapes and detecting skin tone [7,8], will provide greater opportunities for us to explore different kinds of hand gestures in conjunction with verbal commands.

Finally, we have not addressed how context and related discourse issues can affect human-robot interactions in a command and control situation. We are currently adding sensitivity to context into our interface. For example, commands can be interrupted momentarily for some reason, and then resumption of the previous activity may be warranted.

Currently, we are not able to perform such actions. Once we terminate an action, it cannot be resumed, unless that action is repeated. We would like, therefore, to add this capability. Given context, we will be able to issue a command, such as “Go over there,” momentarily stop it, and then resume it, either by issuing a verbal command, such as “Continue,” or simply by offering the appropriate natural gesture which signifies the same thing.

5. CONCLUSIONS

Our goal has been to develop a natural language and gestural interface to a semi-autonomous robot. The use of natural language and gesture in the interface is based on two assumptions. The first is that while natural language is ambiguous, gestures disambiguate certain information in speech. Secondly, humans use natural gestures more easily when giving directive and locomotive commands to a mobile robot. Our interface does not require the user to wear any special gear, other than a headset with which to issue verbal commands, nor to learn a series of symbolic gestures in order to interact visually with the robot. While our corpus of commands is limited, it has been constrained to meet the limitations of the current vision system of the mobile robot.

In the future, we hope to expand the research by incorporating a more robust speech recognition system. This will enable us to increase the types of verbal information that can be input for commands to a mobile robot. We also wish to employ a more sophisticated vision system, such as one capable of detecting the human hand so that more complex gestures can be incorporated with the speech signal. Finally, we continue to expand on the natural language capabilities of the system; namely, we are currently adding context awareness and discourse capabilities to the natural language component.

The authors would like to thank Stephanie Everett, Elaine Marsh, Ken Wauchope, and the researchers of NCARAI for their numerous suggestions in the preparation of this report, as well as having to watch us put Coyote through its paces.

6. REFERENCES

- [1] Kortenkamp, D., Huber, E., and Bonasso, R.P. "Recognizing and Interpreting Gestures on a Mobile Robot," in Proceedings of the Thirteenth National AAAI Conference on Artificial Intelligence, 1996, pp. 915-921.
- [2] Cassell, J. "Speech, Action and Gestures as Context for On-going Task-oriented Talk," in Working Notes of the 1995 AAAI Fall Symposium on Embodied Language and Action, 1995, pp. 20-25.
- [3] Cassell, J., *et al.* "Modeling the Interaction between Speech and Gesture," in Proceedings of the Sixteenth Annual Conference of the Cognitive Science Society, 1994, pp. 153-158.
- [4] Yamauchi, B., Schultz, A., Adams, W., Graves, K., Grefenstette, J., and Perzanowski, D. "ARIEL: Autonomous Robot for Integrated Exploration and Localization," in Proceedings of the Fourteenth National AAAI Conference on Artificial Intelligence and Ninth Innovative Applications of Artificial Intelligence Conference, 1997, pp. 804-805.
- [5] Wauchope, K. "Eucalyptus: Integrating Natural Language Input with a Graphical User Interface," NRL Technical Report, NRL/FR/5510--94-9711, February 25, 1994, Navy Center for Applied Research in Artificial Intelligence, Naval Research Laboratory, Washington, DC.
- [6] Wauchope, K., Everett, S., Perzanowski, D., and Marsh, E. "Natural Language in Four Spatial Interfaces," in Proceedings of the Fifth Conference on Applied Natural Language Processing, Association for Computational Linguistics, 1997, pp. 8-11.
- [7] Erenshteyn, R., Laskov, P., Foulds, R., Messing, L., and Stern, G. "Recognition Approach to Gesture Language Understanding," in Proceedings of the Thirteenth IEEE International Conference on Pattern Recognition, 1996, pp. 431-435.
- [8] Saxe, D. and Foulds, R. "Toward Robust Skin Identification in Video Images," in Proceedings of the Second IEEE International Conference on Automatic Face and Gesture Recognition, 1996, pp. 379-384.